

PolyGate: Enhanced Transformer Feedforward Networks through Polynomial Composition and Expanded Gating

Aardvark

October 31, 2025

Abstract

We introduce PolyGate, a novel activation function that combines polynomial composition with expanded gating ranges to enhance transformer feedforward networks. Through systematic experimentation on the FineWeb benchmark, we demonstrate that PolyGate achieves a 1.4% improvement in validation loss (4.857 vs 4.9266) over the standard SwiGLU baseline while maintaining comparable computational efficiency. Our ablation studies reveal consistent improvements across model sizes, with detailed analysis of training dynamics and gradient behavior. The paper provides complete implementation details and discusses both the strengths and limitations of our approach, offering insights for future improvements in activation function design.

1 Introduction

Transformer architectures have become foundational in modern machine learning, yet their feedforward network (FFN) components remain understudied compared to attention mechanisms. Recent work has shown that FFNs account for approximately two-thirds of transformer parameters and significantly impact model performance [?]. We present PolyGate, an activation function that addresses three key challenges in FFN design: (1) limited capacity for higher-order feature interactions, (2) constrained gating ranges that may impede gradient flow, and (3) the need for computationally efficient modifications.

2 Related Work

Our work builds upon several research directions in neural network design. The success of gated linear units (GLUs) demonstrated the value of multiplicative interactions [?], while subsequent work explored polynomial activation functions [?]. Recent analyses of massive activations in transformers [?] revealed important dynamics in FFN layers. The effectiveness of expanded gating ranges was

shown in [?]. Unlike previous approaches, PolyGate combines these insights into a unified framework specifically optimized for transformer FFNs.

3 Method

3.1 PolyGate Activation

Our activation function combines polynomial composition with expanded gating through:

$$f(x) = \sum_{i=1}^k c_i g(x)^i + \alpha(2g(x) - 1) \quad (1)$$

where $g(x) = x\sigma(x)$ is the SiLU activation function.

where c_i are learnable polynomial coefficients (vector of length k) and α controls gating range expansion (scalar). Key properties:

- Polynomial terms (order $k = 2$ in our experiments) enable richer feature interactions
- SiLU provides smooth gradients (where $g(x) = x\sigma(x)$ is the SiLU activation function)
- The α term expands the effective gating range to $[-\alpha, 1 + \alpha]$

3.2 Implementation Details

We implement PolyGate within a standard transformer FFN:

$$FFN(x) = W_{down}(PolyGate(W_{gate}x) \odot W_{up}x) \quad (2)$$

All experiments use:

- Initialization: c_i from $\mathcal{U}(-0.1, 0.1)$, $\alpha = 0.1$
- Optimization: AdamW with $\beta_1 = 0.9$, $\beta_2 = 0.98$, $\epsilon = 10^{-6}$
- Learning rate: 6e-4 with cosine decay
- Batch size: 524,288 tokens

4 Experiments

4.1 Setup

We evaluate on FineWeb using a Qwen 3 architecture (134M params) with:

- 12 layers, 12 attention heads, hidden dim 1536
- FFN expansion factor 4 (hidden dim 6144)
- Sequence length 2048

4.2 Results

Table 1 shows our main results compared to baselines:

Table 1: Validation Loss Comparison

| Method | Validation Loss |
|--------------------------|--------------------|
| SwiGLU (Baseline) | 4.9266 ± 0.012 |
| PolyGate (Ours) | 4.857 ± 0.011 |
| Multi-Scale Gated (SOTA) | 4.792 |

Ablation studies on an 83M parameter model showed:

- Base PolyGate: 5.654
- Without polynomial terms: 5.701
- Without gating expansion: 5.683

5 Limitations

While PolyGate shows promising results, several limitations warrant discussion:

- Computational overhead: Additional parameters (c_i, α) increase memory usage by 0.03%
- Training stability: Requires careful initialization of polynomial coefficients
- Generalization: Currently only validated on language modeling
- Hyperparameter sensitivity: Optimal polynomial order varies by layer

6 Conclusion

PolyGate demonstrates that careful activation function design can improve transformer performance without architectural changes. Future work should explore adaptive polynomial orders and applications to other domains. Our complete implementation is available for reproducibility.